

Enabling Language Models to Implicitly Learn Self-Improvement

Ziqi Wang, Le Hou, Tianjian Lu, Yuexin Wu, Yunxuan Li,
Hongkun Yu, Heng Ji



Motivation: Enabling models to self-improve

Preference data contains self-improvement signals!

Preference Data:

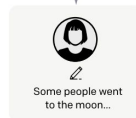
1. Contains a **reference** response and an **improved** response
2. Implicitly contains the **self-improvement** information

Step 1
Collect demonstration data,
and train a supervised policy.

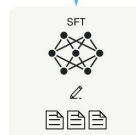
A prompt is
sampled from our
prompt dataset.



A labeler
demonstrates the
desired output
behavior.

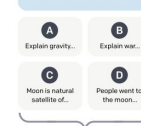
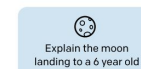


This data is used
to fine-tune GPT-3
with supervised
learning.



Step 2
Collect comparison data,
and train a reward model.

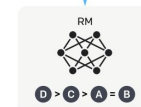
A prompt and
several model
outputs are
sampled.



A labeler ranks
the outputs from
best to worst.



This data is used
to train our
reward model.

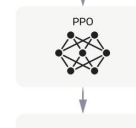


Step 3
Optimize a policy against
the reward model using
reinforcement learning.

A new prompt
is sampled from
the dataset.



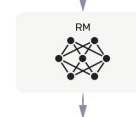
The policy
generates
an output.



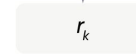
The policy
generates
an output.



The reward model
calculates a
reward for
the output.



The reward is
used to update
the policy
using PPO.



Motivation: Enabling models to self-improve

Preference data can not only be used to train LLMs that generate good responses, but can also train models that generate better responses based on reference responses

$$y_P \sim \mathbf{M}_P(\cdot|x) \longrightarrow y_{\text{PIT}} \sim \mathbf{M}_{\text{PIT}}(\cdot|x, y_{\text{ref}})$$

Why not prompting for self-improvement?

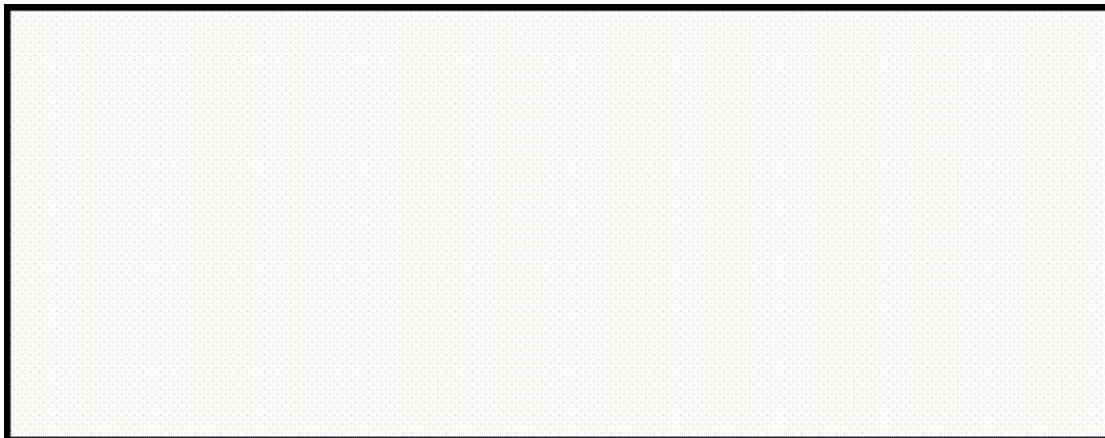
Prompting: “Can you reflect your response and **improve it?**”

Improve:

- More helpful?
- More formal?
- More specific?
- More objective?

More helpful?

- More polite
- More details
- Avoid irrelevant info
- Avoid answering dangerous questions
-



Why not prompting for self-improvement?

Human: What is the origin of the name "Black Friday"?

PaLM 2: ... The first stores to hold such sales in the 1950s called them "Black Friday" because of the increased profits they experienced.

Human: Can you reflect your response and give a more **helpful** response to the question?

PaLM 2:Shopping on Black Friday can be stressful and uncomfortable for many people due to the crowds and crowds, and can also result in overspending and impulse purchasing. For many consumers, the stress of Black Friday shopping outweighs the potential financial benefits, and so shopping on Black Friday is avoided by many consumers.

Why not prompting for self-improvement?

It is generally hard to describe the self-improvement requirements precisely in language.

Learning from preference data is an easier way since models can extract the self-improvement requirements implicitly from data.

Methods

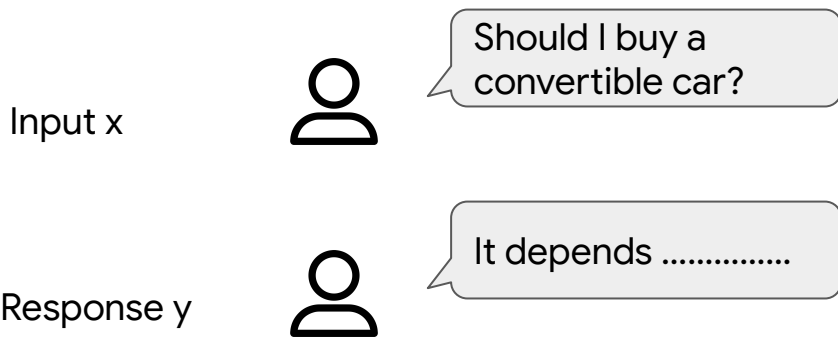
$$y_P \sim \mathbf{M}_P(\cdot|x) \longrightarrow y_{\text{PIT}} \sim \mathbf{M}_{\text{PIT}}(\cdot|x, y_{\text{ref}})$$

Reformulate RLHF pipeline:

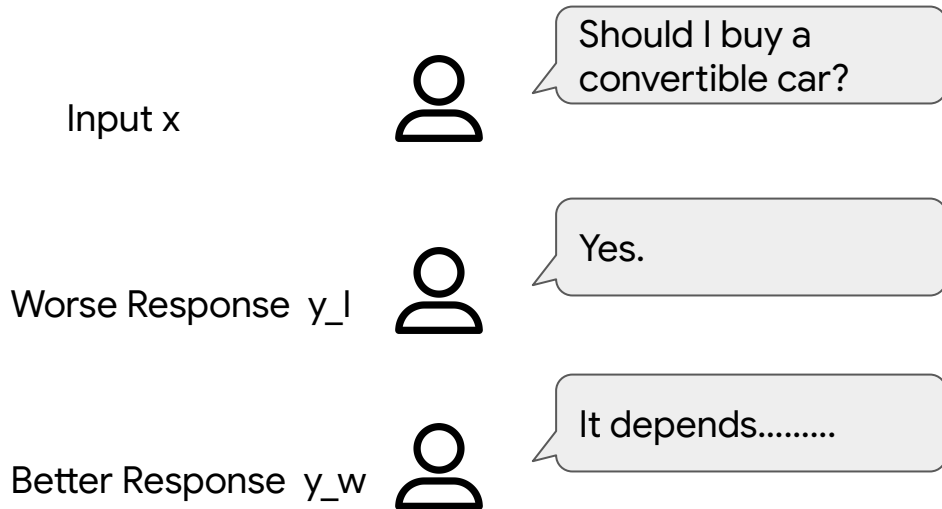
- Supervised Fine-Tuning
- Reward Model Training
- Reinforcement Learning
- Inference

Method

→ **Supervised Fine-tuning:** Learn to generate **<better>** human written responses



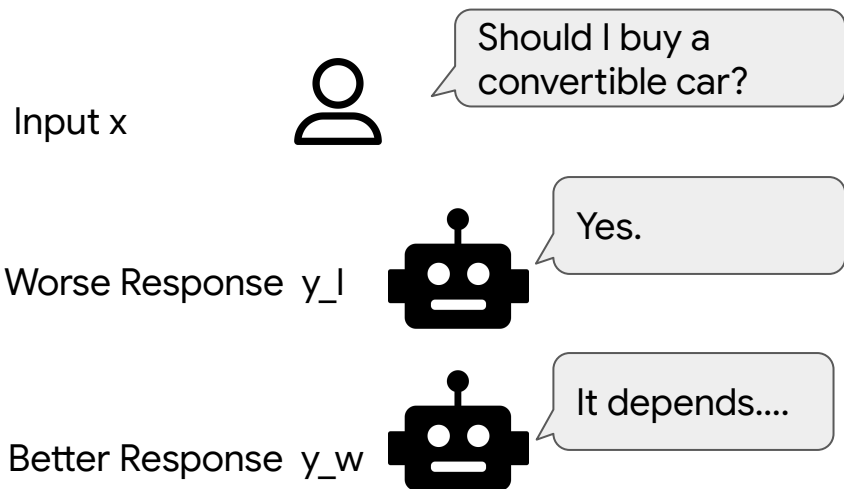
$$-\sum_{(x, y_l, y_w) \in \mathcal{D}_{\text{SFT}}} \log M_P(y_w | x)$$



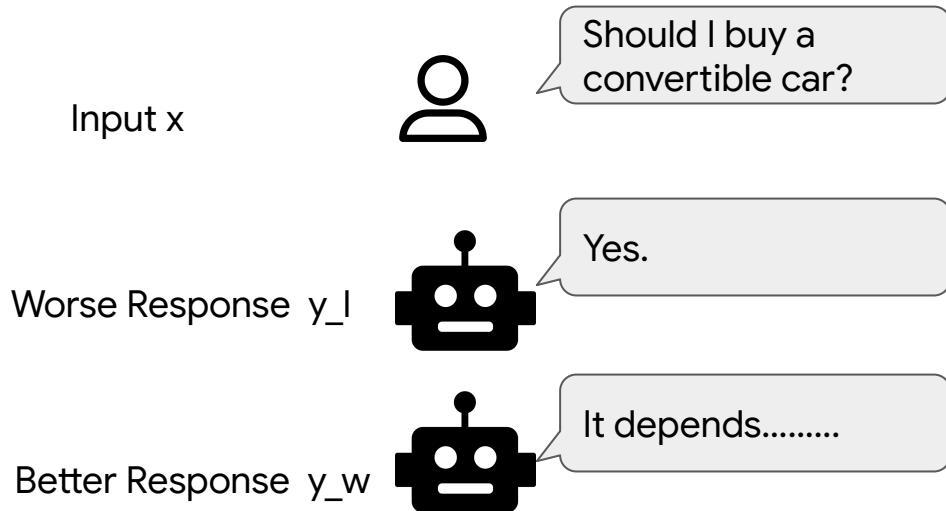
$$-\sum_{(x, y_l, y_w) \in \mathcal{D}_{\text{SFT}}} \log \bar{M}_{\text{PIT}}(y_w | x, y_l)$$

Method

→ **Reward Model Training:** Learn to distinguish between **<better>** and **<worse>** responses



$$-\sum_{\mathcal{D}_{\text{RM}}} \log \sigma(r_w - r_l)$$



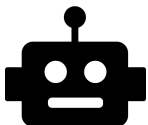
$$r_{\text{gap}}(x, y_w, y_l) \geq r_{\text{gap}}(x, y_w, y_w) \\ \approx r_{\text{gap}}(x, y_l, y_l) \geq r_{\text{gap}}(x, y_l, y_w)$$

Method

→ Reinforcement Learning: Learn to generate **<better>** responses



LLMs



<sample 1>

<sample 2>

<sample 3>

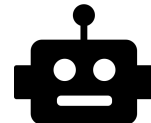
$$\sum_{\mathcal{D}_{RL}} [r(x, y) - \beta \text{KL}(\mathbf{M}_P^{\text{RL}}(y|x) - \mathbf{M}_P^{\text{SFT}}(y|x))]$$

Input x



Should I buy a convertible car?

PIT



<sample 1'>

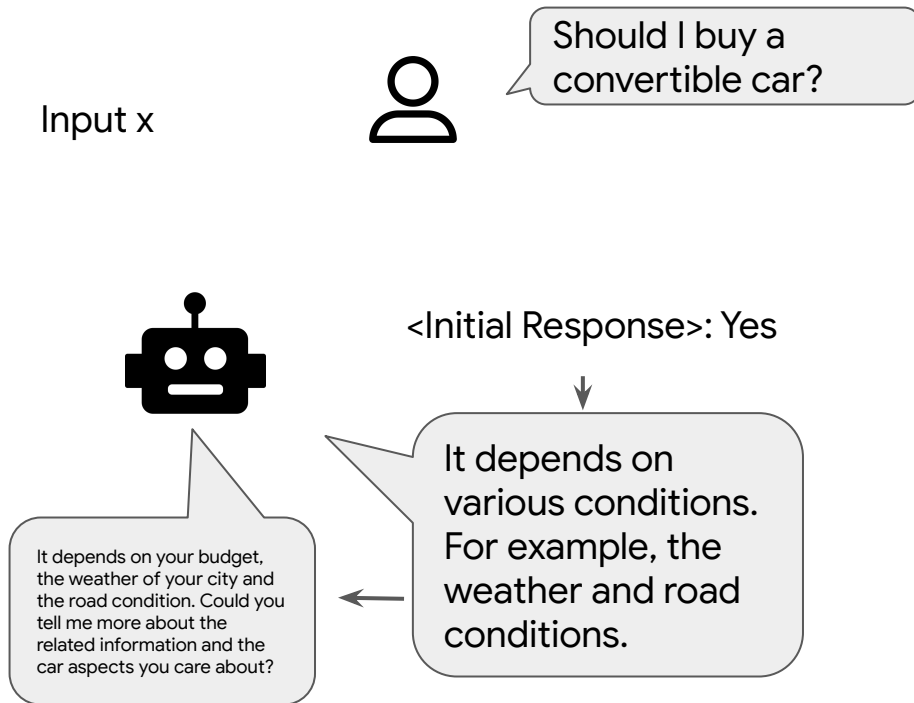
ref

<Initial>

$$\sum_{\mathcal{D}_{RL}} \sum_{y_{\text{ref}} \in \{y_l, y_w\}} [r_{\text{gap}}(x, y, y_{\text{ref}}) - \beta \text{KL}(\mathbf{M}_{\text{PIT}}^{\text{RL}}(y|x, y_{\text{ref}}) - \mathbf{M}_{\text{PIT}}^{\text{SFT}}(y|x, y_{\text{ref}}))]$$

Method

→ Inference: Generate a **<better>** response



Experiment: PIT outperforms Self-Refine

1. PIT improves the LLM response.
2. Self-Refine improves the LLM response, too.
3. PIT performs better than Self-Refine.

Dataset	Comparison	Win rate / Lose rate / Δ (%)		
		GPT-4	DeBERTa	Human Evaluation
Anthropic/HH-RLHF	Original vs. y_w	71.85/17.19/54.69	68.20/18.00/50.20	-
	PIT vs. Original	55.47/27.34/28.13	46.30/32.30/14.00	-
	Self-Refine vs. Original	60.94/17.19/43.75	40.30/31.40/8.90	-
	PIT vs. Self-Refine	38.28/42.19/-3.91	41.3/37.60/3.70	47.06/23.53/23.53
OpenAI/Summary	Original vs. y_w	74.22/8.59/65.63	84.90/10.70/74.20	-
	PIT vs. Original	44.53/24.22/20.31	41.9/34.7/7.2	-
Synthetic Data	Original vs. y_w	28.91/51.56/-22.66	-	-
	PIT vs. Original	48.44/14.84/33.59	-	-
	Self-Refine vs. Original	34.38/17.97/16.41	-	-
	PIT vs. Self-Refine	45.31/35.16/10.16	-	-

Conclusion

$$y_P \sim \mathbf{M}_P(\cdot|x) \longrightarrow y_{\text{PIT}} \sim \mathbf{M}_{\text{PIT}}(\cdot|x, y_{\text{ref}})$$

Thank you!

Paper: <https://arxiv.org/pdf/2310.00898.pdf>

